



EP3260: Machine Learning Over Networks
Computer Assignment 3
Due Date: March 7, 2023

Computer Assignment 3 - Training a neural network

Consider optimization problem

$$\underset{\mathbf{W}_1, \mathbf{W}_2, \mathbf{w}_3}{\text{minimize}} \quad \frac{1}{N} \sum_{i \in [N]} \|\mathbf{w}_3 \mathbf{s}(\mathbf{W}_2 \mathbf{s}(\mathbf{W}_1 \mathbf{x}_i) - \mathbf{y}_i)\|_2^2,$$

where $\mathbf{s}(\mathbf{x}) = 1/(1 + \exp(-\mathbf{x}))$. You may add your choice of regularizer. Using the “Individual household electric power consumption” and “Greenhouse Gas Observing Network” datasets, address the following questions:

- (a) Try to solve this optimization task with proper choices of size of decision variables (matrix \mathbf{W}_1 , matrix \mathbf{W}_2 , and vector \mathbf{w}_3) using GD, perturbed GD, SGD, SVRG, and block coordinate descent. For the SGD method, you may use the mini-batch version.
- (b) Compare these solvers in terms complexity of hyper-parameter tuning, convergence time, convergence rate (in terms of # outer-loop iterations), and memory requirement

• Adding regularizer :

$$\underset{\mathbf{W}_1, \mathbf{W}_2, \mathbf{w}_3}{\text{minimize}} \quad \frac{1}{N} \sum_{i \in [N]} \|\mathbf{w}_3 \mathbf{s}(\mathbf{W}_2 \mathbf{s}(\mathbf{W}_1 \mathbf{x}_i) - \mathbf{y}_i)\|_2^2 + \underbrace{\lambda (\|\mathbf{W}_1\|_2^2 + \|\mathbf{W}_2\|_2^2 + \|\mathbf{w}_3\|_2^2)}$$

(b) According to the figures :

Hyper-parameter tuning:

GD, PGD, and BCD : only 1 hyper-parameter (learning rate) \rightarrow lowest complexity in terms of hyper-parameter tuning

SGD and SVRG : 2 hyper-parameters (learning rate, mini-batch size)

convergence time :

GD, PGD, and BCD : slower

SGD and SVRG : faster

convergence rate :

GD, PGD, and BCD : slower

SGD and SVRG : faster

SGD and SVRG update the weight matrices W_1 and W_2 and the vector w_3 more frequently

Memory requirement :

GD, PGD, and BCD : lowest memory requirement
(only weight matrices and vector are stored)

SGD and SVRG : highest memory requirement

weight matrices & vector + mini-batch of training samples
(in SGD) and control variates (in SVRG),

$$\text{layer-0} = x = f^0$$

$$\text{layer-1} = s(w_1 x_i) = s(w_1 f^0) = f^{(1)}$$

$$\text{layer-2} = s(w_2 s(w_1 x_i)) = s(w_2 f^{(1)}) = f^{(2)}$$

$$\text{layer-3} = w_3 s(w_2 s(w_1 x_i)) = w_3 f^{(2)} = f^{(3)}$$

$$J = \underbrace{\|f^{(3)} - y_i\|^2}_{\text{error}} \rightarrow \nabla J = 2 \text{ error } \nabla f^{(3)}$$

$$\nabla_{w_3} f^{(3)} = f^{(2)}$$

$$\boxed{\text{layer-3-delta} = 2 \times \text{error} \times f^{(2)}} \left(\nabla \text{with respect to } w_3 \right)$$

$$\boxed{\text{layer-2-delta} =}$$

$$2 \text{ error} \times w_3 \nabla_{f^{(2)}} f^{(3)}$$

$$\text{w.r.t. } w_2$$

$$\text{w.r.t. } w_1$$

$$\boxed{\text{layer-1-delta} =}$$

$$2 \text{ error } w_3 w_2 x_i \nabla_{w_1} f^{(1)}$$